**WE CLAIM:**

1.    A data storage system comprising:

a plurality of storage nodes;

data storage mechanisms implemented in each storage

5    node;

a communication medium linking storage nodes; and

a quantity of data distributed across a selected set
of the storage nodes such that the quantity of data
remains available irrespective of the unavailability of

10   one or more of the storage nodes within the selected set.

2.    The data storage system of claim 1 wherein the
data storage mechanisms on at least two storage nodes
collectively implement a unitary volume of network
storage.

3.    The data storage system of claim 1 wherein the
communication medium comprises:

a public network for receiving access requests for
the data storage system; and

5    a private network enabling communication between
storage nodes.

4.    The data storage system of claim 3 wherein the
public network comprises the Internet.

5.    The data storage system of claim 3 wherein the
private network comprises a virtual private network
implemented over the Internet.

6.    The data storage system of claim 1 further
comprising:

communication processes implemented within each of
the storage nodes operable to exchange state information

5    between at least some of the other data storage nodes.

-23-

7.    The data storage system of claim 1 wherein each of the data storage nodes further comprises data structures configured to store state information about one or more other nodes and the communication links between them.

8.    The data storage system of claim 7 wherein the state information comprises information selected from the group consisting of but not limited to: availability information, capacity information, quality of service information, performance information, geographical location information, network topological location information.

9.    The data storage system of claim 8 wherein the set of storage nodes is selected by a first of the storage nodes using the state information stored in the first of the storage nodes.

10.    The data storage system of claim 7 wherein the communication processes implement a repetitive peer-to-peer conversation between the set of storage nodes enabling the state information contained in the state information data structures within each individual node to represent a consistent view of the state of the collection of storage nodes.

11.    The data storage system of claim 1 wherein the network comprises:
    a plurality of first level networks, each first level network coupling multiple storage nodes; and
    a second level network coupling at least two of the first level networks.

12.    The data storage system of claim 11 wherein the first level network comprises a connection selected from

-24-

the group consisting of: Ethernet, fast Ethernet, gigabit
Ethernet, Fibre channel, ATM, firewire, Myernet, SCSI,
5   serial, parallel, universal serial bus, and wireless
networks.

13.  The data storage system of claim 1 further
comprising:
     storage management processes executing on one of the
storage nodes to determine state information about each
5   of the set of storage nodes.

14.  The data storage system of claim 1 wherein the
communication medium comprises a secure communication
medium.

15.  The data storage system of claim 1 wherein the
communication medium implements an authentication
protocol between linked storage nodes.

16.  The data storage system of claim 1 wherein the
communication medium implements cryptographic security
between linked storage nodes.

17.  A method of managing data storage in a network
comprising multiple storage nodes, the method comprising
the acts of:
     communicating a storage request to at least one
5   storage node; and
     causing the at least one storage node to implement
the storage request using an arbitrary subset of the
storage nodes.

18.  The method of claim 17 further comprising:
     communicating state information between the multiple
storage nodes; and

-25-

selecting the arbitrary subset of the multiple
5   storage nodes to be used based upon the state
information.

19. The method of claim 17 wherein the act of
implementing the storage request comprises associating
error checking and correcting (ECC) code with storage
request.

20. The method of claim 19 wherein the ECC code is
stored in a single network storage node and the unit of
data is stored in two or more network storage nodes.

21. The method of claim 17 further comprising:
retrieving a stored unit of data specified by the
storage request; and
verifying the correctness of the stored unit of
5   data;
upon detection of an error in the retrieved unit of
data, retrieving the correct unit of data using data
stored in the others of the arbitrary subset of the
multiple storage nodes.

22. The method of claim 17 further comprising:
attempting to retrieve the stored unit of data from
the arbitrary subset of the multiple storage nodes;
detecting unavailability of one or more network
5   storage nodes; and
in response to detected unavailability, retrieving
the correct unit of data using data stored in others of
the arbitrary subset of the multiple storage nodes.

23. The method of claim 22 wherein the
unavailability is caused by failure of one or more of the
network storage nodes.

-26-

24. The method of claim 22 wherein the unavailability is caused by congestion/failure of a network link leading to one or more of the network storage nodes.

25. The method of claim 17 further comprising moving the stored unit of data from one network storage node to another network storage node after the step of storing.

26. The method of claim 17 further comprising:
communicating state information and storage requests amongst the arbitrary subset of the storage nodes; and
encrypting at least some of the information and
5  storage requests before communicating them between storage nodes.

27. The method of claim 17 further comprising:
communicating state information and storage requests amongst the arbitrary subset of the storage nodes; and
authenticating the communication between storage
5  nodes.

28. A data storage system comprising:
a peer-to-peer network of storage devices, each storage device having means for communicating state information with other storage devices, at least one
5  storage device comprising means for receiving storage requests from external entities, and at least one storage device comprising means for causing read and write operations to be performed on others of the storage devices.

29. The system of claim 28 wherein each of the storage devices comprises means for causing read and

-27-

write operations to be performed on others of the storage devices.

30. The system of claim 28 wherein each of the storage devices comprises data structures defined to configure at least two geographically distant ones of the data storage devices as a unitary volume of storage.

31. The system of claim 30 further comprising:

a network coupling to each of the data storage devices; and

5    a storage controller coupled to the network for logically combining the at least two data storage devices into a single logical storage device.

32. A distributed data storage array comprising:

a plurality of network connected storage nodes;

a network interface within each storage node for receiving data and control information from other storage

5    nodes;

a network interface within at least one storage node for receiving data storage access requests from external sources; and

storage management processes within the at least one

10    storage node operable to distribute data storage for a logically contiguous quantity of data across multiple storage nodes.

33. A data storage system implemented on top of a plurality of networked computer systems and a communication network, wherein each of the networked computer systems implements a storage node and comprises:

5    a processor for processing data according to program instructions;

-28-

a network interface coupled to the processor and the network for communicating data with external entities, including other storage nodes, across the network;

10    memory coupled to the processor, the memory comprising storage space configured to store data and instructions used by the processor;

one or more mass storage devices coupled to the processor;

15    a communication process comprising program instructions executing in the storage node and in communication with the network interface to provide an interface to communicate data storage access requests and responses with the external entities;

20    storage management processes comprising program instructions executing in the storage node and responsive to the received data storage access requests and in communication with the network interface to distribute and coordinate data storage operations with external 25 storage nodes.

34. The system of claim 33 wherein the storage management processes include processes that communicate with the external storage nodes to provide fault-tolerant distribution of data across the a plurality of storage 5 nodes.

35. The system of claim 33 wherein the storage management processes include processes for distributing data redundantly to protect against faults that make one or more storage nodes unavailable.

36. The system of claim 33 wherein the storage management processes includes fault recovery processes, wherein the fault recovery processes respond to a fault condition by communicating with at least one of the 5 external storage nodes to make available a set of data

-29-

that would otherwise be unavailable as a result of the fault condition.